

NCBI Mouse Genome Resources

Peter Cooper
National Center for Biotechnology Information

NCBI

Talk Outline

- **About NCBI**
- **Molecular Databases**
 - GenBank
 - RefSeq
- **Web Access**
 - Entrez
 - BLAST
 - Mouse Genome Resources

NCBI

The National Center for Biotechnology Information (NCBI)

- **Created as a part of NLM in 1988**
 - Establish public databases
 - Research in computational biology
 - Develop software tools for sequence analysis
 - Disseminate biomedical information
- **Tools: BLAST(1990), Entrez (1992)**
- **GenBank (1992)**
- **Free MEDLINE (PubMed, 1997)**
- **Human genome (2001)**

NCBI

Molecular Databases

- Primary Databases
 - Original submissions by experimentalists
 - Database staff organize but don't add additional information
 - Example: **GenBank**
- Derivative Databases
 - Human curated
 - compilation and correction of data
 - Example: **SWISS-PROT**, **NCBI RefSeq mRNA**
 - Computationally Derived
 - Example: **UniGene**
 - Combinations
 - Example: **NCBI Genome Assembly**

GenBank: NCBI's Primary Sequence Database

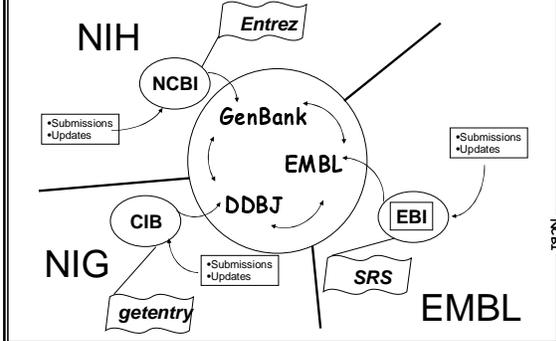
Release 130	June 2002
17,471,130	Records
20,648,748,345	Nucleotides
140,000 +	Species

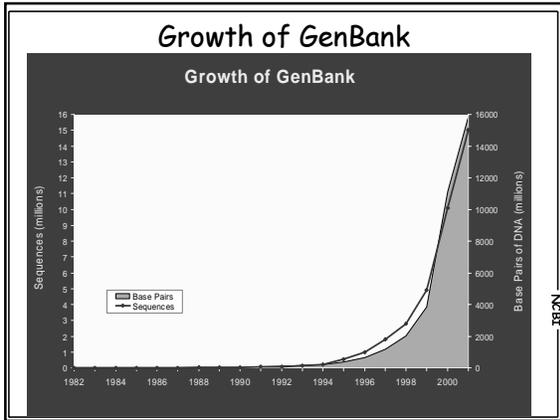
- full release every two months
- incremental and cumulative updates daily
- available only through internet

<ftp://ftp.ncbi.nih.gov/genbank/>

78.75 Gigabytes of data

The International Sequence Database Collaboration





GenBank Divisions

Bulk Sequence Divisions					
PAT	Patent				
EST	Expressed Sequence Tags (142 files)				
STS	Sequence Tagged Sites				
GSS	Genome Survey Sequences (48 files)				
HTG	High Throughput Genome (26 files)				
HTC	High Throughput cDNA				
CON	Contig				
Traditional Divisions					
BCT	INV	MAM	PHG	PLN	PRI
ROD	SYN	UNA	VRL	VRT	

Traditional GenBank Divisions

- Direct Submissions (Sequin and BankIt)
- Accurate
- Well characterized

BCT	Bacterial and Archeal
INV	Invertebrate
MAM	Mammalian (ex. ROD and PRI)
PHG	Phage
PLN	Plant and Fungal
PRI	Primate
ROD	Rodent
SYN	Synthetic (cloning vectors)
VRL	Viral
VRT	Other Vertebrate

Bulk GenBank Divisions

- Batch Submission and htg (email and ftp)
- Inaccurate
- Poorly Characterized

EST	Expressed Sequence Tag
STS	Sequence Tagged Site
GSS	Genome Survey Sequence
HTG	High Throughput Genomic

182N

EST Division: Expressed Sequence Tags

```
>IMAGE:275615 5' mRNA sequence
GACAGCATTGGGCGAGATGTCTCGCTCCGTGGCCCTAGCTGTGCTCGGCTACTCTCTTTCTGGCC
TGGAGGTATCCAGCGTACTCCAAGATTCCAGGTCTACTCACTCATCCAGCAGAGAATGGAAAGTCAAA
TTCTGAAATGCTATGTGTCTGGGTTTCATCCAATCCGACATTGAAGTTGACTTACTGAAGAATGGAGAG
GAATTGAAAAGTGGAGCATTCCAGACTGTCTTTTCAGCAAGGACTGGTCTTTCTATCTTTGTACTACAC
TGAATTCACCCCACTGAAAAGATGAGTATGCCTGCCGTTTGAACCATGTNGACTTTGTTCACAGNCCC
AAGTNGATTTAAGTGGNATCGAGACATGTAAGGCAGGCATCATGGAGGTTTGAAGNATGCCGCTTT
TTGGATGGGATGAATCCCAAATTTCTGGTTTGTCTTGNTTTTTAAATATGGATATGCTTTTG
```

```
-----
>IMAGE:275615 3' mRNA sequence
NNTCAAGTTTATGATTTATTTAATCTGTGGAAACAAAATAAACCCAGATTACCCACAACCATGCCCTACT
TTATCAAATGTAAGANGTAAATATGAATCTTATATGACAAAATGTTTCAATCATTATAACAAAATTC
AATAATCCTGTCAATNATATTTCTAAATTTCCCCCAAATCTAAGCAGAGATGTAATTTGGAAGTTAA
CTTATGACCGCTTAATCTTAAACAAGCTTTGAGTGCAAGAGATTGANGAGTTCAAATCTGACCAAGAT
GTTGATTTGGATAAGGAATTTCTGTCTCCCACTCTANGTTGCCAGCCCTC
```

make cDNA library

80-100,000 unique cDNA clones in library

182N

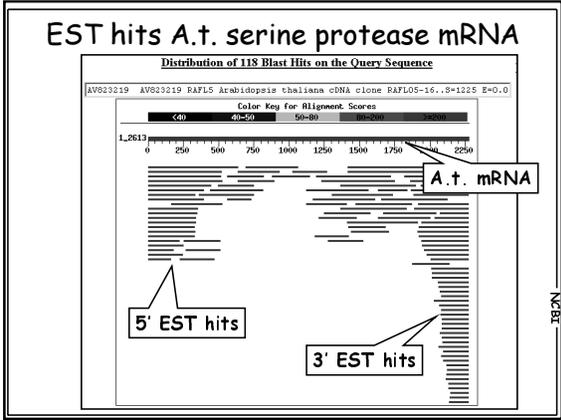
What is UniGene?

A gene-oriented view of sequence entries

- MegaBlast based automated sequence clustering
- Nonredundant set of gene oriented clusters
- Each cluster a unique gene
- Information on tissue types and map locations
- Includes well-characterized genes and novel ESTs
- Useful for gene discovery and selection of mapping reagents

<http://www.ncbi.nlm.nih.gov/UniGene/>

182N



Arabidopsis UniGene Statistics

39,855	mRNAs + gene CDSs	UniGene Build 14
87,006	EST, 3'reads	Apr. 9th, 2002
42,137	EST, 5'reads	
+ 32,571	EST, other/unknown	

201,569	total sequences in clusters	
Final Number of Clusters (sets)		
=====		
26,808	sets total	
	115,000,000 bp	
25,474	25,498 expected genes	one known gene
17,654	5% uncharacterized transcripts	one EST
16,326	sets contain both genes and ESTs	

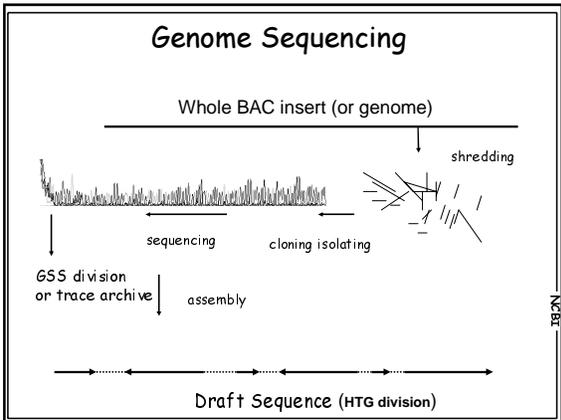
Mm UniGene Statistics

46,745	mRNAs + gene CDSs	UniGene Build 112
1,339,208	EST, 3'reads	Jul. 25th, 2002
89,928	EST, 5'reads	
+ 51,328	EST, other/unknown	

2,332,864	total sequences in clusters	
Final Number of Clusters (sets)		
=====		
84,247	sets total	
	3,000,000 base pairs	
17,723	30 K expected genes	one known gene
82,866	80% uncharacterized transcripts	one EST
16,342	sets contain both genes and ESTs	

UniGene Collections Jul, 2002

		Sequences	Clusters
<i>Homo sapiens</i>	human	3,569,546	101,602
<i>Mus musculus</i>	mouse	2,332, 864	84,247
<i>Rattus norvegicus</i>	rat	334,582	62,220
<i>Danio rerio</i>	zebrafish	197,266	15,404
<i>Bos taurus</i>	cow	128,914	10,295
<i>Xenopus laevis</i>	frog	162,269	18,984
<i>D.melanogaster</i>	fruit fly	250,655	11,115
<i>Anopholes gambiae</i>	mosquito	43,126	2,556
Plants			
<i>Arabidopsis thaliana</i>	thale cress	210,693	26,875
<i>Oryzia sativa</i>	rice	78,632	15,802
<i>Triticum aestivum</i>	wheat	139,447	12,575
<i>Hordeum vulgare</i>	barley	160,518	7,324
<i>Zea mays</i>	maize (corn)	131,668	10,301



HTG Division: High Throughput Genome

LOCUS	AC109609	154774 bp	DNA	linear	HTG 06-FEB-2002
DEFINITION	Mus musculus clone rp23-167h1, WORKING DRAFT SEQUENCE, 19 unordered pieces.				
ACCESSION	AC109609				
VERSION	AC109609.1	GI:18543009			
KEYWORDS	HTG; HTGS PHASE1; HTGS DRAFT.				
LOCUS	AC109609	167074 bp	DNA	linear	HTG 07-MAY-2002
DEFINITION	Mus musculus clone rp23-167h1 strain C57BL/6J, WORKING DRAFT SEQUENCE, 3 ordered pieces.				
ACCESSION	AC109609				
VERSION	AC109609.6	GI:20087170			
KEYWORDS	HTG; HTGS PHASE2; HTGS DRAFT.				
LOCUS	AC109609	167617 bp	DNA	linear	ROD 10-AUG-2002
DEFINITION	Mus musculus chromosome 17 clone rp23-167h1 strain C57BL/6J, complete sequence.				
ACCESSION	AC109609				
VERSION	AC109609.10	GI:22203259			
KEYWORDS	HTG.				

40,000 to > 350,000 bp

Mouse Whole Genome Shotgun

```

LOCUS       CAAA01000000      224713 rc      DNA      linear      ROD 15-MAY-2002
DEFINITION  Mus musculus whole genome shotgun sequencing project.
ACCESSION   CAAA00000000
VERSION     CAAA00000000.1 GI:20800445
KEYWORDS    WGS.
SOURCE      house mouse.
ORGANISM    Mus musculus
             Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
             Mammalia; Eutheria; Rodentia; Sciurognathi; Muridae; Murinae; Mus.
REFERENCE   1
AUTHORS     The Mouse Genome Sequencing Consortium.
TITLE       Direct Submission
JOURNAL     Submitted (01-MAY-2002) The Mouse Genome Sequencing Consortium,
             http://mouse.ensembl.org/
COMMENT     The Mus musculus whole genome shotgun (WGS) project has the project
             accession CAAA00000000. This version of the project (01) has the
             accession number CAAA01000000, and consists of sequences
             CAAA01000001-CAAA01224713.
FEATURES    Location/Qualifiers
             source          1..224713
                       /organism="Mus musculus"
                       /db_xref="taxon:10090"
WGS         CAAA01000001-CAAA01224713
//
    
```

Celera Whole Genome Shotgun

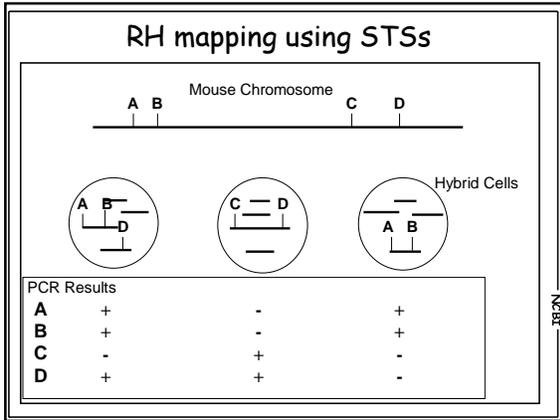
```

LOCUS       AAAD01000000      20 rc      DNA      linear      ROD 31-MAY-2002
DEFINITION  Mus musculus chromosome 16 whole genome shotgun sequencing project.
ACCESSION   AAAD00000000
VERSION     AAAD00000000.1 GI:21281728
KEYWORDS    WGS
SOURCE      house mouse.
ORGANISM    Mus musculus
             Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
             Mammalia; Eutheria; Rodentia; Sciurognathi; Muridae; Murinae; Mus.
REFERENCE   1 (bases 1 to 20)
AUTHORS     Mural,R.J., Chaturvedi,K., Gabriellian,A.E., Ke,Z., Sun,J., Subramanian,G. and
             Venter,J.C.
             Hallberg,S. TITLE A Comparison of Whole-Genome Shotgun-Derived Mouse Chromosome 16
             Sreenivast TITLE and the Human Genome
             Milshina,N. JOURNAL Science 296 (5573), 1661-1671 (2002)
             Wang,K., J. PUBMED 12040188
             Sidick,F. REFERENCE 2 (bases 1 to 20)
             Flanigan,N. AUTHORS Mural,R.J., Adams,M.D., Myers,E.W., Smith,H.O. and Venter,C.J.
             Reinert,K. TITLE Direct Submission
             Lai,S., L. JOURNAL Submitted (24-MAY-2002) Celera Genomics, 45 W. Gude Dr., Rockville,
             Wang,Z.Y., MD 20850, USA
             Wortman,J. COMMENT The Mus musculus whole genome shotgun (WGS) project has the project
             Johnson,J. accession number AAAD01000000. This version of the project (01) has the
             accession number AAAD01000000, and consists of sequences
             AAAD01000001-AAAD01000002.
             The strings of n's in a record represent gaps between contigs, and
             the length of each string corresponds to the length of the gap.
FEATURES    Location/Qualifiers
             source          1..20
                       /organism="Mus musculus"
                       /db_xref="taxon:10090"
                       /chromosome="16"
WGS         AAAD01000001-AAAD01000002
//
    
```

STS Division : Sequence Tagged Sites

- Segment of gene, EST , mRNA or genomic DNA of known position (microsatellite)
- PCR with STS primers gives one product per genome
- Basis of Radiation Hybrid Mapping
 - UniGene
 - Genome Assembly
- Related resource: Electronic PCR

<http://www.ncbi.nlm.nih.gov/genome/sts/epcr.cgi>



Electronic PCR

NCBI UniSTS

Search UniSTS

Human Genome Resources

UniSTS home Submit

FTP site

Statistics

Related sites e-PCR

e-PCR (old)

Map Viewer

LocusLink

Electronic PCR

Query sequence: gi|20632100|emb|CAA01067543.1|, 13514 bases

Mus musculus whole genome shotgun assembly contig 67542, whole genome shotgun sequence

Site (bases)	Marker	Chr.	Organism
3965..4246	D11Bhm52	11	Mus musculus
12507..12639	D11Seg21	11	Mus musculus

Results for D11Bhm52

UniSTS: 141956

D11Bhm52

Primer information

Forward primer: AGCTCAGAGGTGGTGGTCTAG
 Reverse primer: ATCTCAGGAGTGAACCCAG
 PCR product size: 282 (bp), Mus musculus

Mus musculus

Name: D11Bhm52
 Also known as: MGE7317, MGI:106774

Cross References

LocusLink LocusID: 64820
 Symbol: D11Bhm52
 Description: DNA segment, Chr 11, Boehm 52
 Position: 11 45.0 cM

Mapping information

D11Bhm52 MGCv3 Sequence Chr 11 mv
 Map Position: 78991264-78991545 (bp)

Sequence Maps

D11Bhm52 Sequence Map Chr 11 mv
 Position: 80341795-80342076 (bp)

D11Bhm52 MGI Map Chr 11
 Position: 45.00 (cM)

Genetic Map

Electronic PCR results

Genomic (2)
 AL591131.9 5366 .. 5647 Mouse DNA sequence from clone EP21151A11 on chromosome 11, complete sequence [Mus musculus] (114878 bp)

Y14422.1 4189 .. 4470 Mus musculus DNA for retinal protein 09754 bp

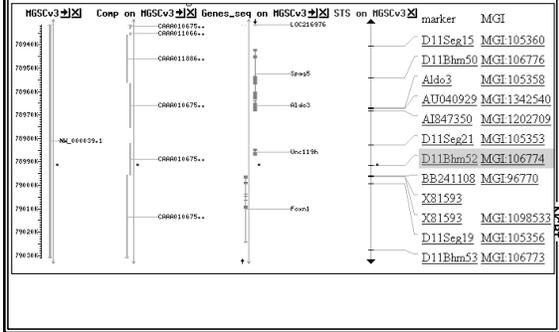
Working Draft phase 1 (from GenBank HTGS division) (2)

AC026375.11 112792 .. 113073 Mus musculus chromosome 11 clone RP23-59385, WORKING DRAFT SEQUENCE, 15 unorderd genes (194695 bp)

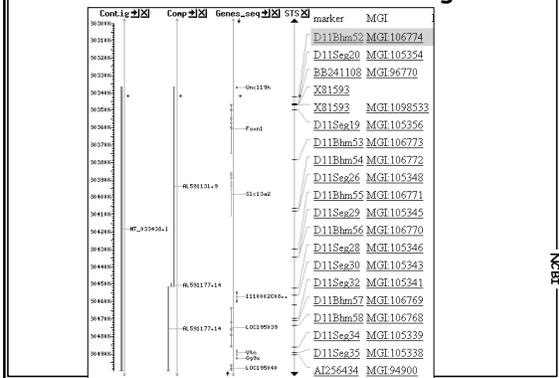
AC004591.1 63640 .. 63921 Mus musculus chromosome 11 clone mCT1.3_M.9 map 11, *** SEQUENCING IN PROGRESS *** 3 unorderd genes (129097 bp)

ePCR Results GenBank Sequences

Markers on MGSCv3

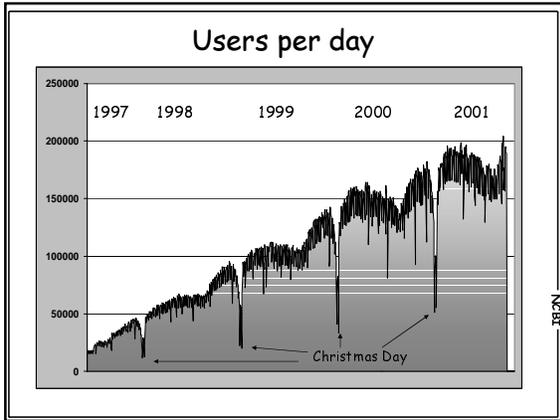


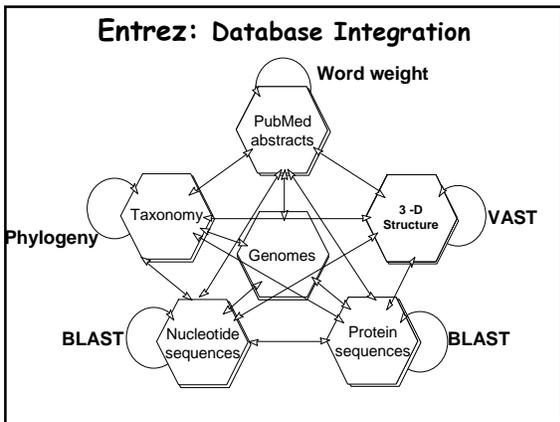
Markers on Finished Contigs



RefSeq: NCBI's Derivative Sequence Database

- **Curated transcripts and proteins**
 - reviewed
 - human, mouse, rat, fruit fly, zebrafish, arabidopsis
- **Human model transcripts and proteins**
- **Assembled Genomic Regions (contigs)**
 - draft human genome
 - mouse genome
- **Chromosome records**
 - microbial
 - organelle





WWW Entrez

The screenshot shows the NCBI Entrez website interface. At the top, there are navigation links for PubMed, Nucleotide, Protein, Genome, Structure, PopSet, Taxonomy, OMM, and Books. Below this is a search bar with a 'Go' button and a 'Clear' button. The main content area is titled 'About Entrez' and 'SITE MAP'. It describes Entrez as a retrieval system for searching several linked databases. A list of available databases is provided, including PubMed, Nucleotide sequence database, Protein sequence database, Structure, Genomes, PopSet, OMM, Taxonomy, Books, and ProbeSet. The interface is clean and organized, with a clear search and navigation structure.

WWW Entrez Data

PubMed: The biomedical literature	All of MEDLINE plus others
Nucleotide sequence database	Abstracts
Protein sequence database	Links to online Journals
Structure: three-dimensional structures	GenBank, EMBL, DDBJ RefSeq, PDB
Genomes: complete genomes	Assemblies
Populations: population statistics	GenBank, DDBJ, EMBL translations PDB, PIR, SWISS-PROT, PRF, RefSeq
OMIM: Online Mendelian Inheritance in Man	
Taxonomy: organisms in the tree of life	Bank
NCBI's MMDB - derived from PDB files	
ProbeSet: gene expression datasets	Reference Genomes: Graphical views, assembled sequence and mapping data
3D Domains: domains from Entrez Structure	
UniSTS: markers and mapping data	
SNP: single nucleotide polymorphisms	
CDD: conserved domains	

Mouse Lipoprotein Lipase

The protein entries in the Entrez search and retrieval system have been compiled from a variety of sources, including SwissProt, PIR, PRF, PDB, and translations from annotated coding regions in GenBank and RefSeq.

Draft Human Genome
Explore human genome resources or browse the human genome sequence using the Map Viewer.

Find Domain Architectures
DART is a way to visualize a protein and its relatives, which takes an accession number or sequence in FASTA format as a query. Hits are returned as schematic diagrams showing the approximate location of domains on the query sequence. Proteins with similar domain architectures are also listed, and can be viewed as subsets - organized by taxonomy or by selected domains.

Retrieve taxonomy information
The Entrez protein database is cross-linked to the Entrez taxonomy database. This allows you to find taxonomy information for the species from which a protein sequence was derived. First, look up a protein in Entrez. A "taxonomy" link appears to the right of each entry that is linked to the Entrez taxonomy database. To view all non-redundant taxonomy links for a search result, select "Taxonomy Links" from the drop-down menu.

Entrez Proteins: Limits

Field Restriction

All fields pull-down menu to specify a field. AND, OR, NOT must be in upper case. Search fields tags are used enclose in square brackets, e.g., rubella [ti].

help on using limits

RefSeq
GenBank
EMBL
DDBJ
PDB
PIR
SWISS-PROT
PIR
PRF

Only from

Date From To

YYYY/MM/DD; month and day are optional.

Entrez Nucleotides: Limits

Limited to:

All Fields Exclude bulk sequences

exclude ESTs exclude STSs exclude GSS exclude working draft
 exclude patents exclude all of the above

Molecule Gene Location Segmented Sequences

Only from Modification Date

Modification Date From To

Use the format YYYY/MM/DD; month and day are optional.

Advanced Search:Preview/Index

Add Term(s) to Query or View Index:

- Enter a term in the text box; use the pull-down menu to specify a search field.
- Click Preview to add terms to the query box and see the number of search results, or click Index to view terms within a field.
- Multiple terms selected from Index will be ORed; click AND to add to search.

Organism Preview Index

Click to add terms selected from Index to the query box.

mouse(125052)

mouse/rat hybrid cell lines(3)

mouse/rat neuroblastoma x glioma hybrid ng108 15(3)

mouse/rat ng108 15(3)

mouse adenovirus(65)

mouse adenovirus 1(65)

mouse adenovirus type 1(65)

mouse cytomegalovirus 1(206)

mouse ear cross(94329)

mouse eared bat(9)

Lipoprotein lipase records

1: AA03305 BLink, Domains, Nucleotide, Related Sequences, Domain Relatives, Taxonomy
 lipoprotein lipase [Mus musculus]
 gi13097036|gb|AA03305.1|[13097036]

2: P11152 BLink, Domains, OMM, Related Sequences, Domain Relatives, PubMed, Taxonomy
 Lipoprotein lipase precursor (LPL)
 gi417255|pp|P11152|LPL_MOUSE[417255]

3: XP_134193 BLink, Domains, Nucleotide, Related Sequences, Domain Relatives, Taxonomy
 lipoprotein lipase [Mus musculus]
 gi20864362|ef|XP_134193.1|[20864362]

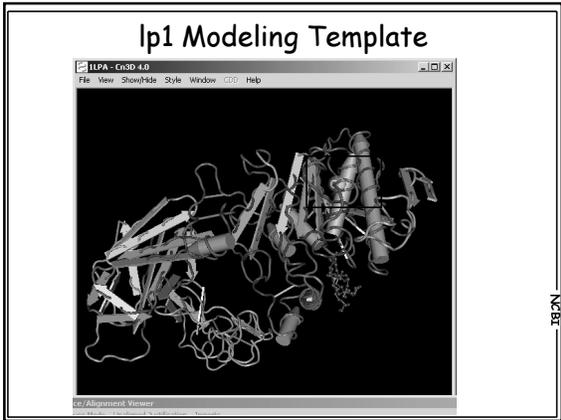
4: NP_032535 BLink, Domains, Nucleotide, OMM, Related Sequences, Domain Relatives, PubMed, Taxonomy
 lipoprotein lipase [Mus musculus]
 gi6678710|ef|NP_032535.1|[6678710]

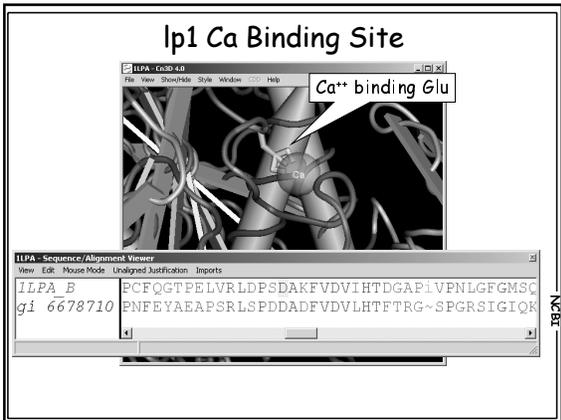
5: A40570 BLink, Domains, OMM, Related Sequences, Domain Relatives, PubMed, Taxonomy
 lipoprotein lipase (EC 3.1.1.34) precursor - mouse
 gi110647|pp|A40570|[110647]

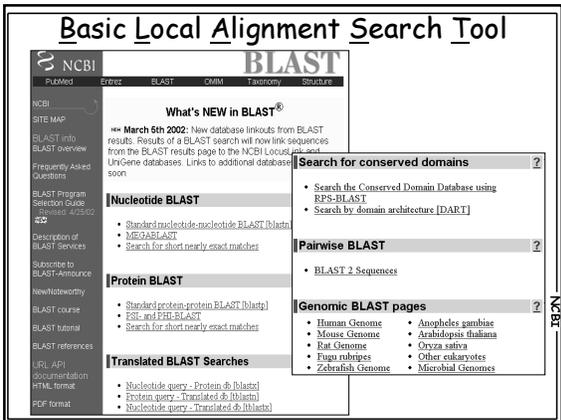
6: S17503 BLink, Domains, Related Sequences, Domain Relatives, PubMed, Taxonomy
 lipoprotein lipase precursor - mouse
 gi110648|pp|S17503|[110648]

7: AAC0464 BLink, Nucleotide, Related Sequences, PubMed, Taxonomy
 lipoprotein lipase [Mus musculus]
 gi2911254|gb|AAC0464.1|[2911254]

8: CAA41329 BLink, Domains, Nucleotide, Related Sequences, Domain Relatives, PubMed, Taxonomy
 lipoprotein lipase [Mus musculus]
 gi51468|emb|CAA41329.1|[51468]







The Standard BLAST Pages

NCBI **protein-protein BLAST**
 Nucleotide Protein Translations Retrieve results for an RID

Search:

Set subsequence From: To:

Choose database:

Do CD-Search:

Now: **BLAST!** or

Advanced Options

Options for advanced blasting

Entrez query restriction: Limit by Entrez query: or select from:

Taxon restriction: Composition-based statistics

Filter options: Choose filter: Low complexity Mask for lookup table only Mask lower case

Expect:

Word Size:

value cut-off: Matrix: Gap Costs: Existence: 11 Extension: 1

BLAST Hits

Taxonomy reports

Distribution of 7 Blast Hits on the Query Sequence

Color Key for Alignment Scores: 0-40, 40-50, 50-60, 60-70, 70-80, 80-90, 90-100

Bit Score: 0, 50, 100

gi|19527038|ref|NP_598579.1|(NM_133818) expressed sequence...135 9e-32

To Entrez	Score (bits)	Significance level
gi 13123094 ref NP_076998.1 (NM_024093) hypothetical prote...	32	1.0
gi 19527038 ref NP_598579.1 (NM_133818) expressed sequence...	31	2.6
gi 18132014 gb AAE73178.1 AF149762.1 (AF149762) decay-accele...	30	4.4
gi 20862594 ref ZP_159074.1 (XM_159074) hypothetical prote...	30	4.7
gi 18132016 gb AAE73179.1 AF149763.1 (AF149763) decay-accele...	29	8.9
gi 20876288 ref ZP_154821.1 (XM_154821) hypothetical prote...	29	8.9
gi 20985709 ref ZP_142189.1 (XM_142189) similar to d310708...	29	8.9

